

This document is published in:

Computer Networks 4 (2010) 12, pp. 2071-2085

DOI: 10.1016/j.comnet.2010.03.022

A collaborative P2P scheme for NAT Traversal Server discovery based on topological information

Rubén Cuevas^{a,*}, Ángel Cuevas^a, Albert Cabellos-Aparicio^b, Loránd Jakab^b, Carmen Guerrero^a

^aUniversidad Carlos III de Madrid, Departamento de Ingeniería Telemática, Spain

^bUniversitat Politècnica de Catalunya, D. d'Arquitectura de Computadors, Spain

Abstract: In the current Internet picture more than 70% of the hosts are located behind Network Address Translators (NATs). This is not a problem for the client/server paradigm. However, the Internet has evolved, and nowadays the largest portion of the traffic is due to peer-to-peer (p2p) applications. This scenario presents an important challenge: two hosts behind NATs (NATed hosts) cannot establish direct communications. The easiest way to solve this problem is by using a third entity, called Relay, that forwards the traffic between the NATed hosts. Although many efforts have been devoted to avoid the use of Relays, they are still needed in many situations. Hence, the selection of a suitable Relay becomes critical to many p2p applications. In this paper, we propose the *Gradual Proximity Algorithm (GPA)*: a simple algorithm that guarantees the selection of a topologically close-by Relay. We present a measurement-based analysis, showing that the GPA minimizes both the delay of the relayed communication and the transit traffic generated by the Relay, being a QoS-aware and ISP-friendly solution. Furthermore, the paper presents the Peer-to-Peer NAT Traversal Architecture (P2P-NTA), which is a global, distributed and collaborative solution, based on the GPA. This architecture addresses the Relay discovery/selection problem. We have performed large-scale simulations based on real measurements, which validate our proposal. The results demonstrate that the P2P-NTA performs similarly to direct communications with reasonably large deployments of p2p applications. In fact, only 5% of the communications experience an extra delay that may degrade the QoS due to the use of Relays. Furthermore, the amount of extra transit traffic generated is only 6%. We also show that the P2P-NTA largely outperforms other proposals, where the QoS degradation affects up to more than 50% of the communications, and the extra traffic generated goes beyond 80%.

Keywords: p2p, Relay, NAT, AS, Simulation, Measurements

1. Introduction

The rapid reduction of available free IPv4 address space [1] has stimulated the widespread deployment of Network Address Translators (NATs) [2] on the Internet. Furthermore, companies' network security policies include the utilization of NATs and/or firewalls in order to hide the

network topology and control both inbound and outbound traffic. As a result, recent studies state that more than 73% of Internet end-hosts are located behind NATs or firewalls [3]. In the following we will use the term *NATed* to designate entities behind NAT.

NATs were designed for client/server applications. However, in the last decade peer-to-peer (p2p) applications such as VoIP (e.g., Skype), online games, P2P file sharing (e.g., BitTorrent), P2P streaming (e.g., PPLive) have become tremendously popular and nowadays they are responsible for the largest share of Internet traffic [4]. Unfortunately, p2p communications cannot be directly

* Corresponding author. Tel.: +34 667718559.

E-mail addresses: rcuevas@it.uc3m.es (R. Cuevas), acrumin@it.uc3m.es (Á. Cuevas), acabello@ac.upc.edu (A. Cabellos-Aparicio), ljakab@ac.upc.edu (L. Jakab), carmen.guerrero@uc3m.es (C. Guerrero).

established through NATs. This is because NATs do not allow inbound connections unless they are manually configured to do so. Researchers have designed a set of techniques to provide NAT traversal capabilities to the NATed hosts [5–10]. However, there are a large number of cases (e.g. symmetric NAT) [11] where these techniques do not work. In these cases the communication must be established using a third non-NATed entity that we call Relay (also known as NAT-Traversal Server). In such a scenario, both end-hosts communicate through this Relay ($end-host1 \leftrightarrow Relay \leftrightarrow end-host2$), that forwards the traffic between them.

As a consequence, the selection of an appropriate Relay becomes a key issue that has a direct impact on the communication delay and, in addition, it may avoid extra-costs for the ISP that hosts the Relay. In particular: (i) a bad choice of the Relay may increase the delay, leading to an undesirable QoS experienced in delay-sensitive communications such as multimedia (VoIP, online gaming, etc.) and (ii) if the Relay is not carefully selected, it increases the total transit traffic of the ISP that hosts it. It must be noted that some p2p applications with a large number of NATed users, such as BitTorrent or eMule, generate a large amount of traffic that may produce an increase of the ISP's costs. Therefore, it is critical to define a Relay selection algorithm that eliminates, or at least alleviates, these negative effects.

This paper presents the *Gradual Proximity Algorithm* (GPA). GPA is a Relay selection algorithm that chooses a topologically close-by Relay from the available pool: first, it tries to find a Relay in the same Autonomous System (AS) as the NATed client. In case there is no Relay in that AS, a Relay within the same NATed node's country is selected; if this fails too, the GPA selects a Relay in the NATed node's continent (to avoid intercontinental links); finally if all the previous attempts fail, a random Relay is chosen. We rely on real measurements to demonstrate that the GPA minimizes the delay of the relayed communication. Also, this algorithm minimizes the extra transit traffic of the Relays' ISPs.

Additionally, this paper presents the Peer-to-Peer NAT-Traversal Architecture (P2P-NTA). This is a wide area collaborative solution that addresses the problem of Relay discovery and selection by using the GPA. Basically, the P2P-NTA is a distributed solution where Relays form a Distributed Hash Table (DHT) to store their location information: AS, country and continent. The NATed nodes are connected to a Relay that belongs to the DHT. In the case when the node is connected to a distant Relay (e.g., on a different continent), the former asks the latter for a closer Relay. The Relay then runs the GPA: (i) it queries the DHT to find a Relay in the same AS as the client; (ii) if there is none, it queries for a Relay in the same country; and (iii) if no Relay is found in the same country, it asks for a Relay on the same continent. After that, the Relays found are sent to the NATed host that connects to the closest one with enough available bandwidth.

Therefore, if we assume that the DHT is well populated and there is at least one Relay per AS, the relayed communications would experience a latency similar to the direct counterpart. This is because we are just adding an additional intra-AS hop that is likely to have a low associated delay. Moreover, all relayed communications follow the

same AS-Path as the direct one, producing no extra transit costs to the ISPs.

In order to validate the P2P-NTA architecture we have developed the P2P-NTA simulator that uses the real Internet AS-topology and real end-to-end latencies. This data has been obtained from the iPlane project [12,13]. In particular, this project provides AS connectivity, IP prefixes announced in the default-free zone and delay measurements between pairs of Points of Presence (PoPs) in the Internet. This dataset has allowed us to run very large-scale simulations involving thousands of real ISPs and several thousands of end-hosts. The results assess the validity of our proposal: its performance is similar to that of direct communication when considering a reasonably sized deployment. Less than 5% of communications suffer from an extra delay that may affect the QoS. Moreover, the P2P-NTA generates an almost negligible amount of transit traffic (6% of the worst-case maximum cost), thus confirming that it is an ISP-friendly solution. On the other hand, the P2P-NTA clearly outperforms other Relay selection algorithms such as Random or Pre-Established Relay selection. In these proposals, more than 50% of the communications suffer from an extra delay that may affect the QoS, whereas even in the best case over 85% of extra transit traffic is generated.

The proposed solution is based on the collaboration of nodes participating in p2p applications (e.g. PPLive, BitTorrent) where many of these end-hosts may act as Relays. Typically, any member of the p2p network with a public IP address is a potential Relay candidate. This leads to a large number of potential Relays; thus, becoming a Relay implies low cost.

In short, the main contributions of this work are:

- *The Gradual Proximity Algorithm*: This is a lightweight and simple algorithm that allows finding a topologically close-by Relay from the available pool. This allows to minimize the relayed communication delay and avoids to generate extra transit traffic at the Relay's ISP.
- *The Peer-to-Peer NAT Traversal Architecture (P2P-NTA)*: This is a globally distributed and collaborative architecture that solves the problem of Relay discovery and selection. For this purpose it uses a DHT to register and retrieve the location information of the Relays, and implements the GPA's selection procedure. This lightweight architecture inherits the advantages of the GPA.
- *The P2P-NTA simulator*: This is an Internet-scale simulator that allows us to evaluate the architecture considering thousands of real ISPs and tens of thousands of final users. Since it uses real AS-maps and real end-to-end delays the results are relevant. It must be noted that this simulator can be adapted to evaluate other large-scale solutions in the Internet, so we release it to the scientific community.¹

The rest of the paper is organized as follows. Section 2 revises the related work. In Section 3, we discuss why it

¹ The simulator is available at <http://personals.ac.upc.edu/acabello/p2pnat>.

is smart to select a topologically close-by Relay and present our *Gradual Proximity Algorithm* for Relay selection. Section 4 describes the Peer-to-Peer NAT-Traversal Architecture. We devote Section 5 to detail the trace-based simulator used for the evaluation of our solution, whereas Section 6 shows the results of the evaluation. Finally, Section 7 concludes the paper.

2. Related work

2.1. Graph representation and network coordinate systems (NCSs)

The networking research community has dedicated some effort to predict the Internet graph. These works give an estimation of the distance between hosts in the Internet. A first approach consists of creating a real Internet graph. Two of the main representatives of this approach are IDMaps [14] and iPlane [12]. The former uses an infrastructure formed by a set of Vantage Points distributed around the Internet, named *Tracers*. The tracers measure the distance among them. The rest of the hosts are clustered in reachable Address Prefixes (APs). Furthermore, the system measures the distance between each AP and its nearest tracer. Hence, the distance between two APs can be calculated as the sum of the distance from each AP to its nearest tracer plus the distance between the tracers. iPlane [12] is based on a similar concept. However, the system implements sophisticated measurement techniques to estimate the delay, loss rate, capacity and bandwidth of the path between two Internet end-hosts. These systems need a dedicated measurement infrastructure and incur a high probing load.

On the other hand, the NCSs do not need an infrastructure of Vantage Points to perform measurements. In this approach, the probing is performed by all the nodes involved in the system. The aim of the NCSs is to map the system (e.g. Internet) topology into a multi-dimensional coordinate system where each node has associated, given virtual coordinates. For this purpose, each node performs measurements to other nodes in the systems in order to find its correct position. Based on the virtual coordinates each node can estimate the distance to any other node in the virtual coordinate space. Vivaldi [15] is the most known representative of this family. More recent works have improved Vivaldi and applied the NCSs to the Azureus DHT [16] and online games applications [17]. Although the NCSs do not need a dedicated measurement infrastructure they still cause a high probing load, furthermore they are not robust against the Triangular Inequality Violation (TIV) and tend to fall into unoptimum local minimum states.

These solutions help to identify the location of a given node as well as identify the distance between nodes. Therefore, it would be feasible to apply them to the problem of Relay Selection. However, they are more complex and cause a much higher probing load than our proposed solution. Additionally, they are ISP-unaware, thus they must be redesigned in order to be ISP-friendly.

2.2. Overlay routing

It is commonly accepted by the research community that with the current AS-based routing (BGP [18]), the direct route between two end-hosts may be suboptimal in terms of delay. Hence, several solutions have been proposed, that may reduce this delay in some cases, using an overlay routing approach. These solutions use one or multiple intermediate overlay relay nodes in order to shorten this AS-Path. Solutions such as RON (where the Relays are selected from a static pool) [19] or SORS (where the Relays are selected at random) [20] are intended to improve general IP routing. Other solutions such as ASAP [21] are designed only for VoIP applications. In particular, ASAP is a complex architecture to find Relays that, according to their results, reduces the delay in some cases. However, it is not clear what the signaling overhead produced by the proposed system is, and it has the disadvantage of relying on bootstrapping servers and cluster representatives (called surrogate nodes) that are single points of failure. Unlike the P2P-NTA, none of these solutions consider the NAT-Traversal problem. In addition, these solutions produce extra transit traffic at the ISPs where the Relays are located, imposing extra costs to these companies.

2.3. Close-by server selection

Guton et al. presented an early work on static and centralized location of nearby servers of a distributed service [22]. The solution combined traceroutes and hop-count measurements to determine the closest replica. One year later, Carter et al. [23] demonstrated that dynamic server selection is more efficient than static server selection due to the variability of route latency over time and the large divergence between hopcount and latency. In parallel, IP Anycast was proposed as a network-layer solution to server selection. It was first proposed in 1993 by the IETF RFC 1546. However, due to various deployment and scalability problems [24], it has not been widely deployed. More recent solutions, Meridian [25] and OASIS [26], also address this problem. In Meridian, the servers form an Overlay where each server knows some other servers and locate them in concentric rings based on the measured RTT. When a given client launches a query to find a close-by server, the query progresses through the overlay until it reaches the closest server. During the process a large number of servers have to measure the RTT to the client, which constitutes a high probing load. On the other hand, OASIS is an anycast infrastructure issued for multiple services. It has a central infrastructure of nodes used to locate a close-by server to a given client. The delay measurement is performed by the replica servers of the different services registered in OASIS. The measurement procedure is based on an optimization of Meridian.

All the described solutions are designed for classical services where the server is expected to be an *always online* machine. Although they could be applied to p2p applications, the probing load would increase dramatically due to the users churn. Furthermore, in the case of OASIS a dedicated infrastructure of core nodes is needed. Again it is

worth noting that these solutions have been designed neither for ISP-friendliness nor for dealing with NATed hosts.

2.4. Relay/super node selection

To the best of the authors' knowledge there are very few proposals that address the problem of Relay selection in case of hosts behind NATs. The P2PSIP Working Group (WG) [27] of the IETF is currently designing a p2p version of the SIP framework of protocols, where the users are able to establish communications among them without using any rendezvous server such as SIP Proxies or SIP Registrars. For this purpose, they use a DHT. One of their major challenges is how to solve the problem of NAT Traversal. In [28], the authors propose a lightweight mechanism to discover Relays within the P2PSIP DHT. Although the protocol is promising, it does not consider any kind of location information for selecting a Relay, and this incurs the costs already mentioned in this paper. We believe that the P2P-NTA is a good candidate solution to be considered by the P2PSIP WG to solve the problem of NAT-Traversal server discovery.²

Authors in [29] propose VIP, a p2p communication platform for NAT Traversal. In their solution, the nodes use ICE [7] and Hole Punching [8–10] in order to traverse the NATs. Basically, the nodes learn their available IP addresses/ports and register them in a DHT, so their *buddies* can easily access this information. Some of these VIP nodes act as Relays for those which are behind NATs. Unfortunately, the paper does not specify how a VIP node discovers one of these Relays.

Finally, some p2p applications such as Skype select as *super nodes* those that show stability and have enough available bandwidth. These *super nodes* act as Relays for that specific application. Skype is a proprietary application, and it is unknown how the Relay is selected. Nevertheless, some researchers have reverse engineered it [21,30,31] and have found that Skype uses Relays even if direct communications are possible. In particular the clients try to establish the connection through different Relays (sometimes dozens of them) before selecting one. It is also known that the Relays are not randomly selected and that the selection is AS-unaware [21,31].

2.5. Locality solutions for P2P applications

Over the last years ISPs are experiencing an extra transit traffic due to p2p applications. Furthermore, some ISPs have started to throttle traffic from some applications such as BitTorrent [32,33]. As a reaction to this problem, recently some works have appeared describing locality solutions to keep the traffic of p2p applications within the local ISP as much as possible [34,35]. To the best of our knowledge, there is no previous work that addresses this issue for relayed communications. Thus, we believe that our proposal is the first contribution regarding transit traffic reduction for relayed communications in the Internet.

² Due to lack of space we cannot explain the details of how to specifically implement P2P-NTA within the P2PSIP architecture.

3. Smart Relay selection: the Gradual Proximity Algorithm

In this section, we first discuss why Relays are a must in the current Internet. Next, we clarify why selecting a topologically close-by Relay is efficient and we describe the *Gradual Proximity Algorithm (GPA)*. Finally, we present a measurement-based analysis that validates the proposed algorithm.

3.1. The need of relays

Around 73% of Internet users are located behind NAT [3]. It has been shown that even using very sophisticated techniques [7–10], there are some types of NATs (e.g Symmetric NAT), widely deployed [11], that can be hardly traversed. In addition, not all the applications implement these NAT traversal techniques. This leads to the conclusion that Relays are a necessity of today's Internet.

3.2. Selecting a topologically close-by Relay

When two nodes (*A* and *B*) that are connected behind a NAT want to communicate through a Relay they establish two different connections: $A \leftrightarrow \text{Relay} \leftrightarrow B$. Some previous overlay routing proposals suggested to select the Relay randomly [20], or from a pre-established pool [19]. These selection algorithms may obtain an unsuitable Relay that increases the communication delay and the transit traffic of the ISP that hosts the Relay. Indeed, it is possible that two end-users located in the same ISP choose a Relay from a different one, or even from a different continent.

We claim that selecting a topologically close Relay reduces both the delay of the communications and the transit traffic of the ISP that hosts the Relay. The Internet is structured into Autonomous Systems (AS), hence the closest Relay, in terms of hops, is usually located in the same AS as the node itself. Unfortunately, not all of the ASes may host a Relay. For this case we have defined the *Gradual Proximity selection Algorithm (GPA)*, that defines different degrees of proximity in the current Internet scheme (Algorithm 1).

Algorithm 1. Gradual Proximity Algorithm

```

/* Input Parameters*/
userAS, userCountry, userContinent
/* Output Parameters*/
Relay
/* Algorithm*/
if  $\exists$  Relay on the same AS then
    Choose the closest one with enough bandwidth
    among them
else if  $\exists$  Relay on the same country then
    Choose the closest one with enough bandwidth
    among them
else if  $\exists$  Relay on the same continent then
    Choose the closest one with enough bandwidth
    among them
else
    Choose a random Relay with enough bandwidth
end if
return Relay

```

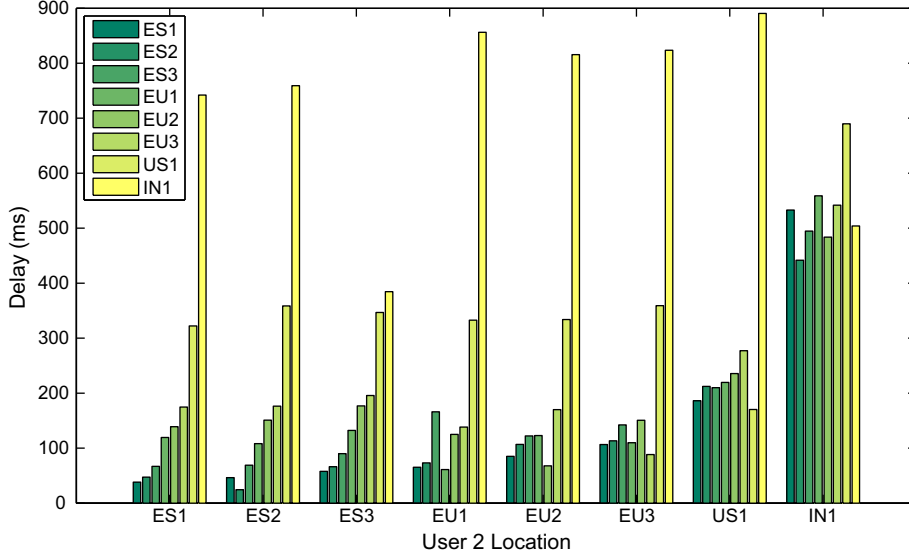


Fig. 1. RTT of relayed communication ($User1 \leftrightarrow Relay \leftrightarrow User2$). User 1 is always located in ES1. x-axis identifies User 2 AS. Each bar represents the Relay's AS as indicated by the legend.

Next, we explain the consequences that the GPA has in terms of delay and transit traffic cost.

3.3. Consequences of the GPA in terms of delay and transit traffic cost

The idea behind our algorithm is to use the shortest possible AS-Path between the two end hosts, the main rationale being the following: (i) if we use a Relay in the same AS, the AS-Path will be the same as for direct communication. Then, GPA adds just some extra hops inside the AS at IP level; (ii) if we choose a Relay in the same country we are probably adding one AS-hop toward the Relay³; and (iii) if we select a Relay on the same continent, it is likely to add some hops to the AS-Path, but we avoid intercontinental links (e.g., transoceanic) that have a very high propagation delay.

To validate the performance of the proposed algorithm in terms of delay, we have performed a set of live experiments. We have deployed measurement boxes in different ASes: 3 in Spain: ES1, ES2, ES3; 3 in three different European countries: EU1 (Italy), EU2 (Norway), EU3 (Greece); 1 of them located in the US, US1; and the last one located in India, IN1. We use the notation $u1$ and $u2$ for the end users involved in the communication and R for the Relay. The location of $u1$ is fixed to ES1 while we iterate $u2$ and R among all the possible locations. Fig. 1 shows the RTT values of the relayed communications between $u1$ (always located in ES1) and $u2$ (located in the AS indicated by x-axis) through R (located in the AS indicated by the legend).⁴

³ It is likely that two ISPs within the same country have a peering agreement, thus being one AS-hop apart to each other.

⁴ We measured the RTT, for each end-user to end-user communication, 10 times per day at different day hours during one week. Fig. 1 shows the average value of the RTT.

As the figure shows, the topological distance between the Relay and the end-users has a significant impact on the communication delay. For instance, IN1 at the largest topological distance from any ES or EU location. Thus, for all those cases where $u1$ and $u2$ are located anywhere in Europe, using a Relay in India produces the highest delay. We can also see that, for a given $u2$ location, the latency increases as assumed by the GPA: same AS < same country < same continent < different continent. The results also suggest that there is a positive relation between the topological distance and the delay.

To further validate our algorithm we have measured the end-to-end delay of a large set of end-users. The measurements come from iPlane [12,13], which is a scalable service, providing accurate predictions of Internet path performance for overlay services and Internet-scale simulations. To achieve these goals, the iPlane project uses hundreds of vantage points distributed across the Internet for measurements, updating their dataset daily. iPlane is based on daily active latency measurements from various vantage points of the Internet. In particular they take advantage of the PlanetLab infrastructure and, using trace-route, they monitor hundreds of paths from each of the available PlanetLab nodes. In this experiment we have obtained a latency dataset by querying the iPlane service using random IP addresses. The iPlane interface includes in each reply a flag indicating if the requested latency has been either estimated or measured. In our dataset we only consider measured latencies.

In particular, we have measured the one-way delay for: (i) end-users located in the same AS; (ii) end-users located in the same country but different ASes; (iii) end-users located on the same continent but different country; and (iv) end-users located on different continents. Our dataset contains 1M end-to-end delays. Fig. 2 shows the Empirical Cumulative Distribution Function (ECDF) of the delay for the different cases. In this context, the term delay refers

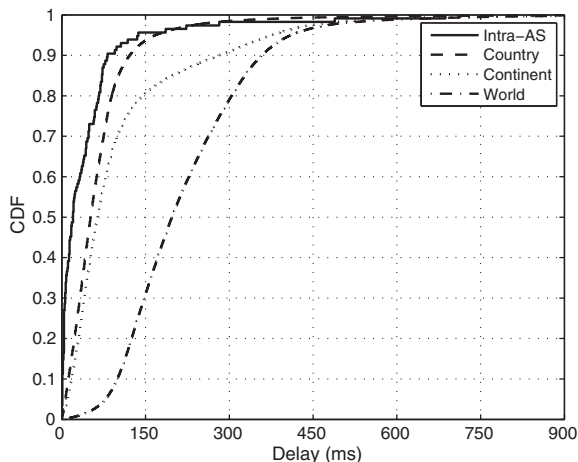


Fig. 2. One way delay for different topological end-user locations.

to the latency measured by iPlane using traceroute, and targets the instantaneous delay.

These distributions suggest that the delay is strongly related to the topological distance between the end-users. If we consider the maximum one-way delay recommended by the ITU-T for voice communications (150 ms) [36], we can conclude according to this experiment that more than 95% of the communications are below that threshold, if both end users belong to the same AS. This percentage is still over 90% when the users are located in the same country but in different ASes. When the users are located on the same continent, the percentage drops to 83%. Finally the delay is severely impacted (40%) if the users are on different continents, and the end-to-end path includes trans-continental links.

It is also worth to note that while ITU recommends 150ms as a quality threshold for voice communications, empirical experiments run by Cisco show that there is a negligible difference in voice quality Mean Opinion Score (MOS) when a 200 ms threshold is used [37]. Throughout the paper we will use a 200 ms delay threshold to differentiate between acceptable and non-acceptable quality voice communications. In particular, in this dataset if we consider the Cisco's threshold: 98%, 96%, 87%, and 50% of the communications are below 200 ms when users are located in the same AS, same country, same continent and different continents, respectively.

Therefore, the obtained results confirm the suitability of GPA as Relay selection algorithm that effectively minimizes the end-to-end communication delay.

3.4. Consequences of GPA for transit traffic

In this subsection, we evaluate the extra traffic caused by the GPA algorithm. ISPs usually pay both for inbound and outbound traffic flowing through transit links, while the cost of traffic transmitted through peering links is typically free [38]. In order to better understand how the GPA avoids the transit traffic, let us focus on the case where $u1$ and $u2$ are located in different countries and ASes. In this

case, what is the extra cost that the Relay's ISP has to pay compared to the case of direct communication?:

- (1) If we choose a Relay in a country and AS other than the two end-hosts, the Relay uses transit links to communicate with both users. Therefore, the ISP where the Relay is located has to pay for the out-bound and inbound traffic of $u1$ and $u2$.
- (2) If we select a Relay in a different AS, but in the same country as $u1$, the Relay uses a peering link to communicate with $u1$ ⁵ whereas it uses a transit link to communicate with $u2$. Thus, the cost for the ISP where the Relay is located is half than in the previous case.
- (3) If we select a Relay in the same AS as either $u1$ or $u2$, the relayed AS-Path is the same as the direct one, thus the ISP does not incur any extra cost.

Fig. 3 shows the different scenarios explained above. As a result, the Gradual Proximity Algorithm always selects a Relay that minimizes the ISP transit traffic, and therefore can be considered as an ISP-friendly algorithm.

In a nutshell, we have demonstrated that our GPA leads to a reduction in the communication delays while minimizing the transit traffic costs compared to other proposals.

Finally, it must be highlighted that these benefits could have a direct impact in currently deployed applications with millions of users. On the one hand, VoIP applications such as Skype use real-time communications that are delay-sensitive, therefore GPA would improve the quality of the communications. On the other hand, the majority of the users of P2P file sharing applications such as eMule or BitTorrent are located behind NATs,⁶ thus requiring a Relay. Moreover, these applications produce a large amount of traffic, and this increases the costs for ISPs. In this scenario, the GPA reduces significantly the relayed transit traffic and the costs for ISPs with regard to other proposals.

4. The P2P NAT-Traversal Architecture (P2P-NTA)

In this section, we present the *P2P Nat Traversal Architecture (P2P-NTA)*. This is a collaborative, distributed and wide-area application that implements the Gradual Proximity Algorithm as described in Section 3. First, we detail the proposed architecture and its functionality.

4.1. Physical architecture

The *P2P-NTA* consists of two different types of nodes, *clients* and *Relays*. Nodes located behind NAT are clients, while nodes having a public IP address and enough available bandwidth typically become Relays. The latter form a Distributed Hash Table (DHT) where they register their location information. Each client has an associated Relay from this DHT for NAT-Traversal capabilities. Fig. 4 depicts the physical architecture.

⁵ Access ISPs within the same country typically establish peering agreements.

⁶ We have conducted a large-scale crawling of BitTorrent, demonstrating that more than half of the users are located behind NATs. More detailed information can be found in [39].

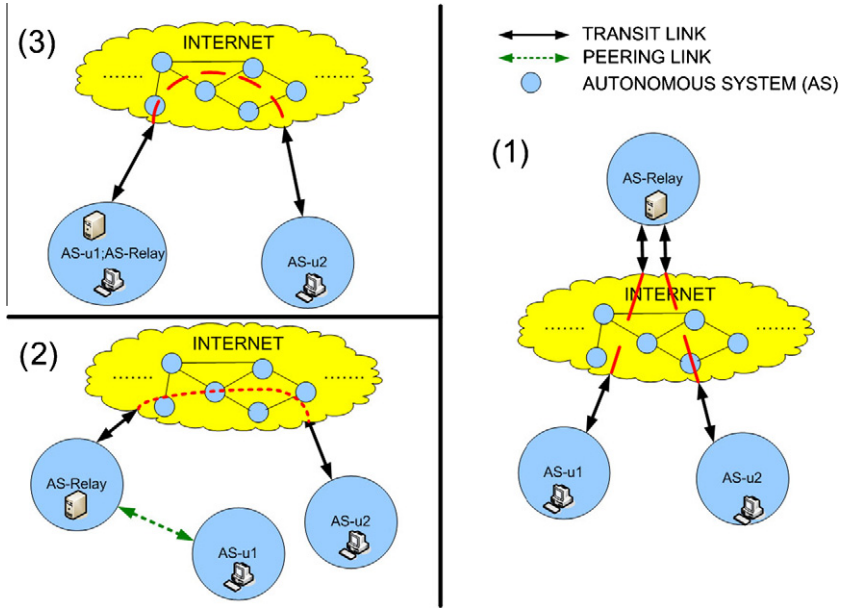


Fig. 3. ISP-friendly Relay selection. *Case 1:* the selected Relay is located in a different AS and country as the users. It uses (paid) transit links to both users. *Case 2:* the selected Relay is located in an AS in the same country as one of the users. It uses a (paid) transit link to one of the users and a peering (free) link to the other. *Case 3:* the Relay is located in the same AS as one of the users. The relayed communication follows the same AS-Path as the direct communication, thus incurring in no extra cost.

4.2. Bootstrapping

When turned on, the node checks if it is a Relay (it has a public IP address) or a client (it is behind a NAT). If it is a Relay it joins the DHT. For this purpose, it computes its *Peer-ID* as $\text{hash}(\text{IP address})$. This *Peer-ID* indicates the position of the node in the DHT. After that, the Relay contacts

any member of the DHT (previously known, well-known peers, Bootstrapping server, etc.) to join the *P2P-NTA*. On the other hand, if the node is a client it needs to attach to a Relay belonging to the *P2P-NTA*. If it is the first time the node joins the P2P, it will contact a bootstrapping server or a well-known DHT peer, if not, it will connect to any Relay known in the past. After that, the client checks if this

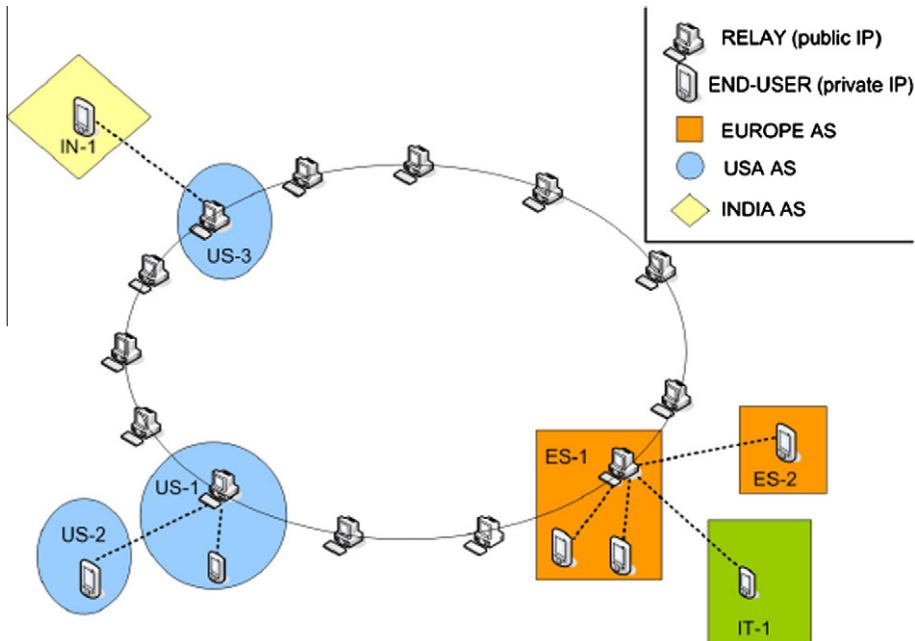


Fig. 4. P2P-NTA physical architecture.

Relay belongs to its own AS and if it provides a delay below a pre-established threshold. In this case, the client keeps this Relay, otherwise it finds a closer one as described below (Relay Lookup Procedure).

4.3. Registration of the Relay's location information

Once the Relay has joined the DHT, it computes its AS key = $\text{hash}(\text{AS number})$, country key = $\text{hash}(\text{country ID})$ and continent key = $\text{hash}(\text{continent ID})$. Then it stores the following tuples: (AS-key, Relay IP), (country-key, Relay IP) and (continent-key, Relay IP) in the DHT. In particular these tuples are stored on the nodes with the closest IDs to the AS-key (AS-key's responsible node), country-key (country-key's responsible node) and continent-key (continent-key's responsible node), respectively.⁷

4.4. Relay lookup procedure (based on GPA)

The client can eventually detect that its current Relay is located in another ISP, or that it is too distant (e.g. the RTT towards the Relay is above a predefined threshold). This triggers the lookup procedure for a closer Relay: the client sends a message to its current Relay, including its publicly visible IP address. In turn, the Relay discovers the client's AS, country and continent and with this information, it computes an AS key, a country Key and a continent Key. Next, the Relay applies the GPA: (i) first, it sends a query to the DHT looking for the AS key. If the query succeeds, the Relay obtains a list with all the IP addresses of the Relays located in the same AS as the client. (ii) If this first query fails, the Relay sends a second query looking for the country Key, if this query succeeds the Relay retrieves the list of IP addresses of Relays located in the same country as the client. (iii) If the query fails again, the Relay sends a third query asking for the continent Key. The Relay, in turn, forwards this information to the client. Thus, the client obtains a list of, either all the Relays located in its AS (if the first query succeeded), all the Relays located in its country (if the second query succeeded) or all the Relays located on its continent (if the third query succeeded). If none of the queries succeed the client keeps its current Relay (instead of selecting a random one as explained in Section 3). It could be possible that the number of Relays located in a given AS, country or continent is too large (hundreds or even thousands). In this case the responsible node would not answer with the full list but rather a limited number (e.g., 50) selected at random.

4.5. Relay selection

The client has to select a Relay from those included on the list. In order to select the best one in terms of available bandwidth and delay, it contacts first the closest one (i.e. the one offering the lowest RTT) and asks to join. If that given Relay has enough available bandwidth, it ac-

cepts the client. Otherwise the client is rejected and tries the second closest Relay, repeating the process until it gets accepted.

4.6. Communication establishment

After the Relay selection algorithm each user is bound to a specific Relay. When two users ($u1$ and $u2$) want to establish a communication they use standard mechanisms, such as ICE [7], to exchange information about their respective Relays ($R1$ and $R2$). With this information the users can test both paths ($u1-R1-u2$ and $u1-R2-u2$) and choose the best Relay, in terms of delay, to establish their connection.

4.7. Geolocation procedure

In our solution, the Relays must be capable of identifying the country, continent and AS number associated to the client's IP address. For this purpose different public services [40] and databases [41,42] can be used.

4.8. Churn and replication

The P2P-NTA users are, in fact, end users that may join and leave the system at any moment. This phenomenon is known as churn. In the P2P-NTA, when a Relay leaves the system gracefully, it notifies the necessary DHT nodes and removes the registered information contacting the nodes responsible of its AS, country and continent. Furthermore, it informs its neighbors and, if necessary, reassigns its stored information (AS, country and continent keys) to them. If the Relay leaves the system abruptly, it leaves the DHT with inconsistent information. To deal with such cases we have defined the following mechanisms:

- **Tuple timer:** The responsible node of an AS, country or continent tuple uses an expiration timer. The Relays must update the tuples before the timer expires, otherwise the responsible node removes the tuple. This way, if a Relay leaves the system abruptly the tuples related to it would automatically expire after a certain amount of time.
- **Replication:** The responsible node replicates the stored information in R replicas to the DHT. These are nodes with the i th ($i \in [2, R + 1]$) closest IDs to the given key. Thus, if the responsible node unexpectedly leaves the system, it is not affected at all, since the first replica takes the responsibility of its keys. In addition, the use of replicas enable load balancing mechanisms [43–46] to share the load among all of them.

Finally, when a Relay leaves the system, its clients must select a new one by triggering the *Relay Selection* algorithm as described above. The worst case happens when a Relay, which is forwarding traffic from two different users, leaves the system abruptly. Since each user is bound to a given Relay, and they have agreed on using one of them during the *Communication Establishment* procedure, they can switch immediately to the other one and resume their communication.

⁷ The *responsible node* refers to the node of the P2P that stores and is authoritative for the requested information.

5. The P2P-NTA simulator

This section describes the P2P-NTA simulator used to validate the proposal. This is an iterative simulator, implementing a Chord DHT [47] with users deployed in the real Internet topology. Moreover, it is using real latencies to account for the communication delays among the nodes.

The foundation of the simulator is the iPlane [12,13] platform (described in Section 3). It builds the topology based on this dataset, which contains AS connectivity, the IP prefixes announced in the BGP default-free zone and delay measurements between pairs of Points of Presence (PoPs) in the Internet. Each PoP, as defined by iPlane, is a set of IP addresses with low latency among them. The simulator uses the PoPs to build the Internet-topology, where we consider 55,000 PoPs and their actual point of attachment. Note that an AS may contain more than one PoP. The P2P-NTA simulator considers these PoPs as the access routers of ISPs and therefore, it deploys the Relays randomly among them. We consider four cases: 100, 1000, 10,000 and 25,000 Relays, each case referring to the amount of Relay nodes contained in the Chord DHT.

Then, for each iteration, the P2P-NTA chooses two different random users from two different PoPs. The users are chosen according to the following criteria, in order: within the same country, within the same continent or from different continents. With this set of experiments we aim to show the performance of our proposal under different scenarios. Nevertheless it is important to remark that the first case, where the users are chosen within the same country, is the most common one, especially in VoIP applications.

The users query the Chord network deployed among the Relays that run the GPA selection algorithm. The P2P-NTA simulator implements a highly scalable Chord network and it is able to route the query towards its destination, and provide the path, number of hops, and latency. In order to simulate very-large Chord networks, the P2P-NTA uses a steady-state approach, and only simulates this P2P network after it has stabilized. We assume that during the simulation no churn is observed, and iteratively generate the finger table for each node in turn, knowing a priori the full list of the nodes in the overlay. After the finger tables are generated, queries can be routed by the simulator using this topology. It is worth noting here that we validated our implementation of the Chord protocol in steady state with OpenChord 1.0.5.⁸ Specifically, we created a P2P network using OpenChord and waited until the network converged. Next, we compared the finger tables of the P2P-NTA simulator and OpenChord nodes, which were found identical.

Finally, once the query has finished, and both users have agreed on communicating using a given Relay with the help of the GPA, the simulator computes both the direct and the relayed delay.⁹ In order to estimate these laten-

cies, iPlane provides the delay between PoPs, but not the delay between the PoP and the end-user (i.e. the access link). This part of the end-to-end delay is estimated using the dataset provided in [48], that measures the median access link speed for different countries. Hence, the user is geolocated [41]¹⁰ and the access link latency is estimated.

Preliminary experiments showed us that, for the scenarios simulated, iPlane provided the latency between two PoPs approximately in 70% of the cases. That is why we included a latency estimator for the remaining 30% cases. In order to design it we have used a dataset that contains roughly 200k latencies¹¹ between arbitrary pairs of hosts. We have divided this dataset (randomly) into two sets, one for training and designing the estimator, and the other one for validation purposes.

In order to design a latency estimator we take into account the information that we can associate to each PoP. In particular we aim to correlate the geographical distance between them with the latency and we consider the following estimators. First a linear regression, secondly we bin the pairs of PoPs depending on their distance and we compute the Empirical Cumulative Distribution Function (ECDF) of the latencies. Then, considering this training data, we estimate the delay of a pair of PoPs firstly computing the distance between them, and then generating a random number that follows the ECDF of the appropriate bin. In particular we consider two bin sizes: (i) (0–10 km, 10–100 km, 100–1000 km, 1000–10,000 km, 10,000–20,000 km) and (ii) (0–10 km, 10–100 km, 100–500 km, 500–1000 km, 1000–2500 km, 2500–5000 km, 5000–7500 km, 7500–10,000 km, 10,000–15,000 km, 15,000–20,000 km). With this approach, we assume that there is a correlation between a given bin and the latency, for instance routers that are at a range of 10 km may have the same amount of hops on their paths. Further, we also consider this approach taking into account the particularities of their location at a continent-level. It is clear that the topology of the Internet is different if we consider North America or Europe, mainly because some continents are more densely populated, and routers may be deployed closer.

Fig. 5 shows the error of the estimators. Each curve represents the ECDF of the absolute error (estimated-real) of the different proposed estimators. The absolute error has been computed subtracting the real latency from the estimated one (from the validation set). As we can see the accuracy of the estimators is similar, except for the *distance/c* estimator, which always under-estimates. This is because it only considers the propagation delay, and assume that end-to-end paths are just a link. Our simulator implements the linear regression estimator since it is the most accurate. Further, this estimator is very fast, and will not slow down the simulator. It is important to note that generating random numbers that follow a certain ECDF is computationally intensive. Also, as Fig. 5 shows, the linear regression estimator is not biased, and since we plan to carry out a large amount of repetitions, this will not impact

⁸ <http://open-chord.sourceforge.net>.

⁹ We define direct delay as the latency between two nodes in the Internet that communicate using the standard inter and intra-domain routing protocols. We define relayed delay as the latency between two nodes that communicate through a third node.

¹⁰ This database is open source and has an 99.8% accuracy at country level, 75% accuracy at city level (within a range of 25 miles), 22% accuracy at more than 25 miles, and 3% that the IP is not covered by such database.

¹¹ A subset of the dataset described in Section 3.

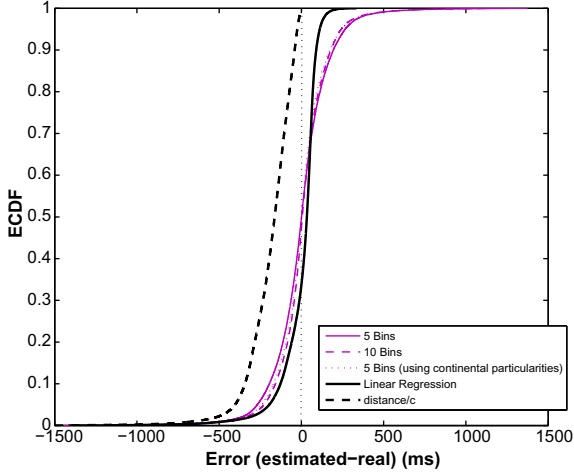


Fig. 5. Error of the different estimators. Each curve represents the CDF of the absolute error of the different proposed estimators.

the results. More details about this estimator and the simulator can be found in [49].

6. Evaluation

In this section, we present the obtained results from our large-scale simulations in terms of: (i) delay; (ii) ISP-friendliness; and (iii) overhead produced by the P2P-NTA (lookup latency and Relay load).

6.1. Simulation set-up

We have used the P2P-NTA simulator described in Section 5 to simulate our proposal with different deployments: 100, 1000, 10,000 and 25,000 Relays. These cases refer to the amount of nodes that can potentially act as Relays because they are configured using a routable IP address and are not firewalled. These nodes, which represent the P2P-NTN service, have been randomly deployed based on the 55,000 real PoPs considered by the iPlane dataset. The clients of the overlay network using the P2P-NTN service are also randomly deploy among the PoPs. Table 1 describes these scenarios in detail, showing how many ASes, countries and continents contain at least one Relay node.

We have also simulated the random Relay selection algorithm for the same number of Relays and the pre-established Relay selection algorithm for a fixed pool of 1000 Relays. These solutions are equivalent to SORS [20] and RON [19], respectively.

We have simulated, for each solution and number of Relays: (i) 30 k communications in the *Intra-Country* scenario: the communication is established between hosts located in the same country; (ii) 30 k communications in the *Intra-Continent* scenario: same continent (but different countries); and (iii) 30 k communications in the *Inter-Continent* scenario: different continents. In total we have simulated roughly 700 k communications to evaluate our solution.

For each communication we calculate the direct and the relayed delay. Also, we geolocate (AS, country and conti-

ment) the two users ($u1$ and $u2$) and the Relay (R) involved in the communication in order to estimate the *ISP-friendliness* for the different solutions. In addition, we calculate the load supported by each Relay in terms of number of relayed communications. Finally, we compute the Relay lookup latency for each communication.

6.2. Communication delay

We have computed the ECDF of the one-way delay for each solution (P2P-NTA, Random Relay Selection and Pre-Established Relay Selection), deployment (100, 1000, 10,000 and 25,000 Relays) and scenario (Intra-Country, Intra-Continent and Inter-Continent). Fig. 6 summarizes the obtained results:

- Fig. 6(a) shows the ECDF for the delay for the 90k communications. We consider the Cisco's 200 ms quality threshold described in [37]. That is, we consider that 200 ms is the maximum tolerable one-way delay for voice communications with acceptable QoS. As expected, the direct communications present the lowest delay, and 79% of them are below the aforementioned threshold. The P2P-NTA slightly increases the direct communication delay under the largest considered deployment. For instance, in the case of 25,000 Relays,¹² 75.5% of the communications are below the 200 ms threshold. This means that less than 4% of the communications would suffer from QoS degradation due to the use of Relays compared to the direct one. As the deployment decreases, the number of communications below the threshold slowly decreases (72.6% for 10,000 Relays and 72.2% for 1000 Relays) up to 56.3% in the 100 Relays case. However, even in such small deployments, the P2P-NTA clearly outperforms other proposed algorithms. The figure shows that Random¹³ and Pre-Established selection algorithms can only keep the communication delay below the 200 ms threshold in 18.3% and 17.2% of the cases, respectively. These values are less than 1/2 of our proposal's worst case (100 Relays).
- Fig. 6(b)–(d), depict the result for the Intra-Country, Intra-Continent and Inter-Continent scenarios, respectively. For clarity we just plot the deployments of 25,000 and 100 Relays for the P2P-NTA solution. The rest of deployments lay between these two curves. In those cases, the amount of communications which are below the quality threshold between the direct and the P2P-NTA (25,000 Relays) is always smaller than 6%. Hence, independently of the type of communication considered (short, medium or long distance) we can conclude that our solution causes a minimum QoS degradation.

Moreover, the P2P-NTA outperforms other Relay selection algorithms. Even if we compare our solution using

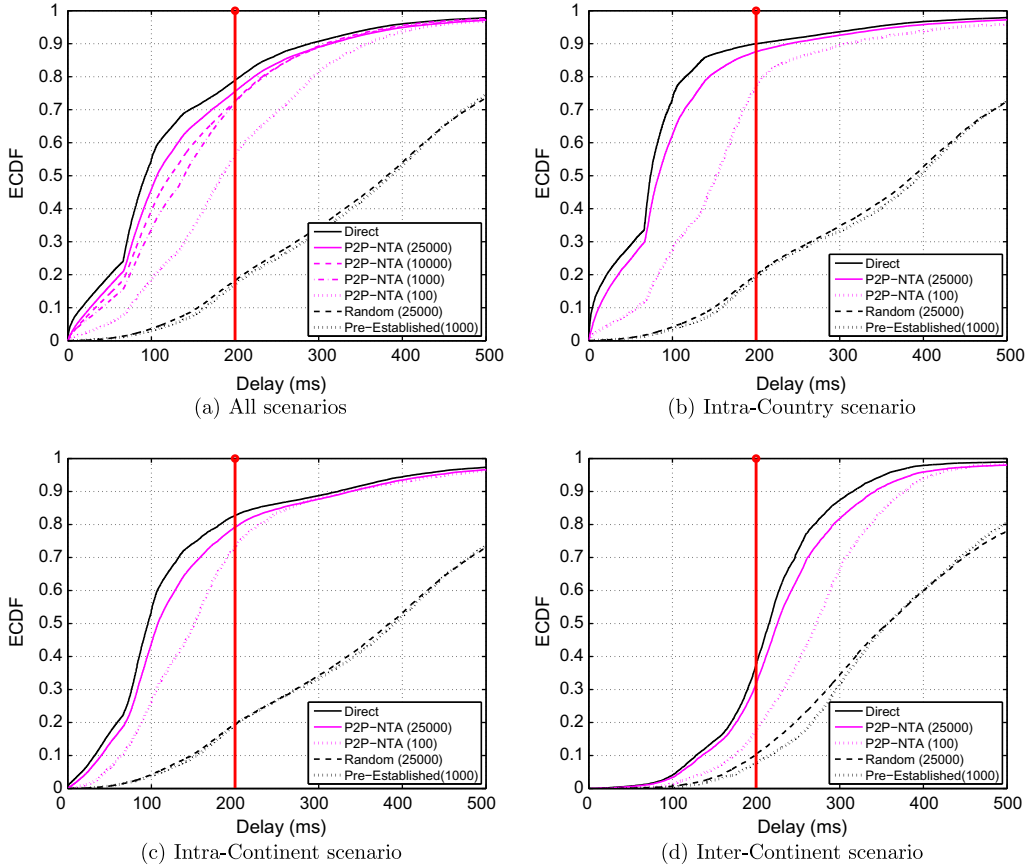
¹² It must be highlighted that 25,000 Relays is actually a small deployment if we consider current p2p applications such as Skype or KAD that include millions of concurrent users.

¹³ Note that the random selection algorithm performs similarly regardless of the number of Relays. Thus, we only depict one case.

Table 1

Distribution of clients and Relays through ASes, countries and continents in the different scenarios (100, 1 K, 10 K and 25 K Relays). Note that the maxmind database considers that South and North America are different continents.

	Clients loc. all scenarios	Relays loc. 100 Relays scenario	Relays loc. 1 K Relays scenario	Relays loc. 10 K Relays scenario	Relays loc. 25 K Relays scenario
#AS	15,815	84	820	4996	10,361
#Countries	203	26	92	149	190
#Continents	6	6	6	6	6

**Fig. 6.** Relayed communications delay.

the smallest deployment (100 Relays), the number of communications below 200 ms is four times larger in the Intra-Country and Intra-Continent case, and almost two times larger in the Inter-Continent case.

Finally, and as a side-result, it is worth to note the problems that long distance delay-sensitive communications (e.g. VoIP) may experience in the current Internet. Fig. 6(d) shows that even in the direct communication case, just 1/3 of the communications are below the 200 ms threshold.

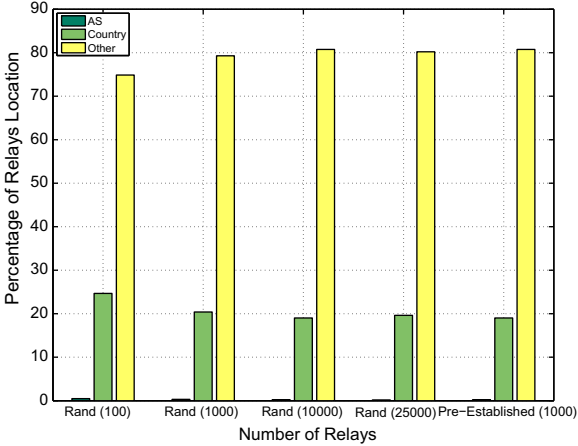
In short, we can conclude that in terms of end-to-end communication delay our solution largely outperforms other previous works, and for (p2p applications) reasonably sized deployments the performance is similar to the direct communications.

6.3. ISP-friendliness

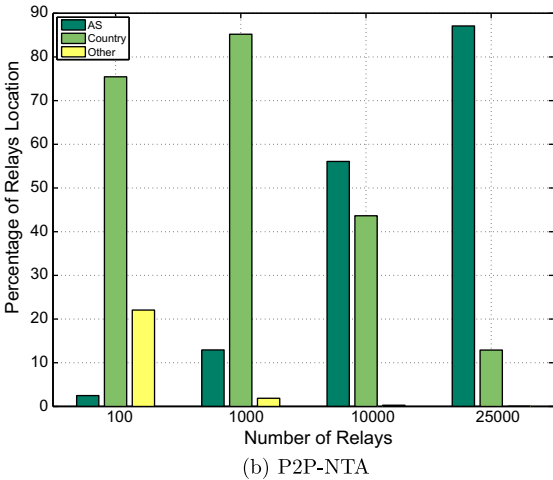
As stated in Section 2, ISPs have recently shown their concerns because of the large amount of traffic generated by p2p applications (e.g., BitTorrent). Fig. 7 supports our initial hypothesis, and shows that the P2P-NTA is ISP-friendly since it minimizes the relayed traffic transmitted through transit links. This figure represents the percentage of communications (out of the 90 k) in which the Relay has been selected in: (A) the AS of one of the two end-users, thus leading to zero cost; (B) an AS in the same country of one of the two hosts, thus using a free peering link to communicate to one user and a paid transit link to communicate to the other; and (C) any other case where the Relay uses transit links to communicate with both end-hosts.

In Fig. 7(a) we present the results for Random and Pre-Established selection algorithms, whereas Fig. 7(b) depicts the results for the different deployments of our solution. We can see that both, Random and Pre-Established selection algorithms, (in the best case) use less than 1% and 25% of Relays located in the same AS (case A) and country (case B) as $u1$ or $u2$, respectively. In other words, they use in all the cases more than 73% of Relays located in different AS and country as both users (case C). This means using transit links towards both $u1$ and $u2$ (case C).

On the other hand, and as expected, the performance of our solution is affected by the deployment. In the case of 25,000 Relays, 87% of the communications are established through a Relay located in the same AS than one of the hosts (case A) while the remaining 13% use a Relay in the same country (case B). If we consider the minimum deployment case with just 100 Relays, our solution still performs quite well since just around 20% of connections are established through a Relay outside the ISP and country of both hosts (case C).



(a) Random and Pre-established Relay Selection



(b) P2P-NTA

Fig. 7. ISP-friendliness. Percentage of selected Relays located in the same AS, same country or different country as the users for the different selection mechanisms.

To further understand the results, we consider relative transit costs: $cost = 0$ for case A; $cost = 1/2$ for case B; $cost = 1$ for case C. Recall that in case B, the Relay uses a transit link to communicate with one of the hosts, whereas in case C it uses transit links for both hosts. Furthermore, case A is free of cost since the Relay is located in the same AS than one of the end-users. Table 2 shows the average cost of the communications for each solution. We observe that Pre-Established and Random selection algorithms are close to the maximum cost (between 87% and 90% of the maximum possible cost), whereas the P2P-NTA produces almost no transit traffic cost (6%) for the largest considered deployment. The cost increases as we reduce the deployment. However, even in the case of minimum deployment (100 Relays), our solution reduces around 30% of the transit traffic compared to the other proposals.

6.4. P2P-NTA overhead

In this section, we evaluate the different aspects of the overhead introduced by our solution. First, we evaluate the Relay lookup latency. Next, we evaluate the number of communications transmitted through the Relay nodes, and we compare them to a Random selection approach.

- *Relay lookup latency:* We define the lookup latency as the time to search through the P2P-NTA and retrieve a list of Relays. We have measured the Relay lookup latency for all the communications (90k) and the different deployments. Fig. 8 summarizes the results. It

Table 2

Average transit traffic cost of the communications (Max = 1, Min = 0).

	100 Relays	1000 Relays	10,000 Relays	25,000 Relays
P2P-NTA	0.60	0.44	0.22	0.06
Random selection	0.87	0.89	0.90	0.90
Pre-established selection	–	0.89	–	–

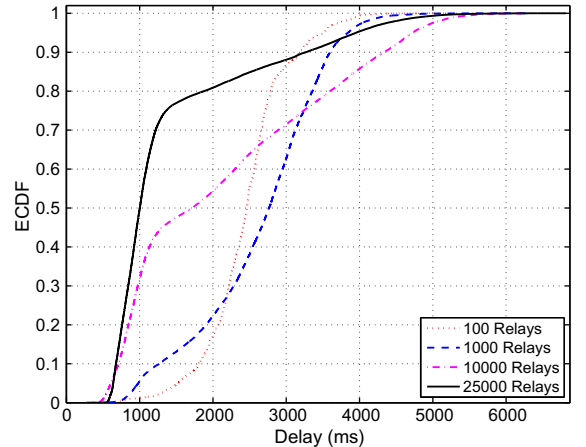


Fig. 8. Relay lookup latency for different deployments.

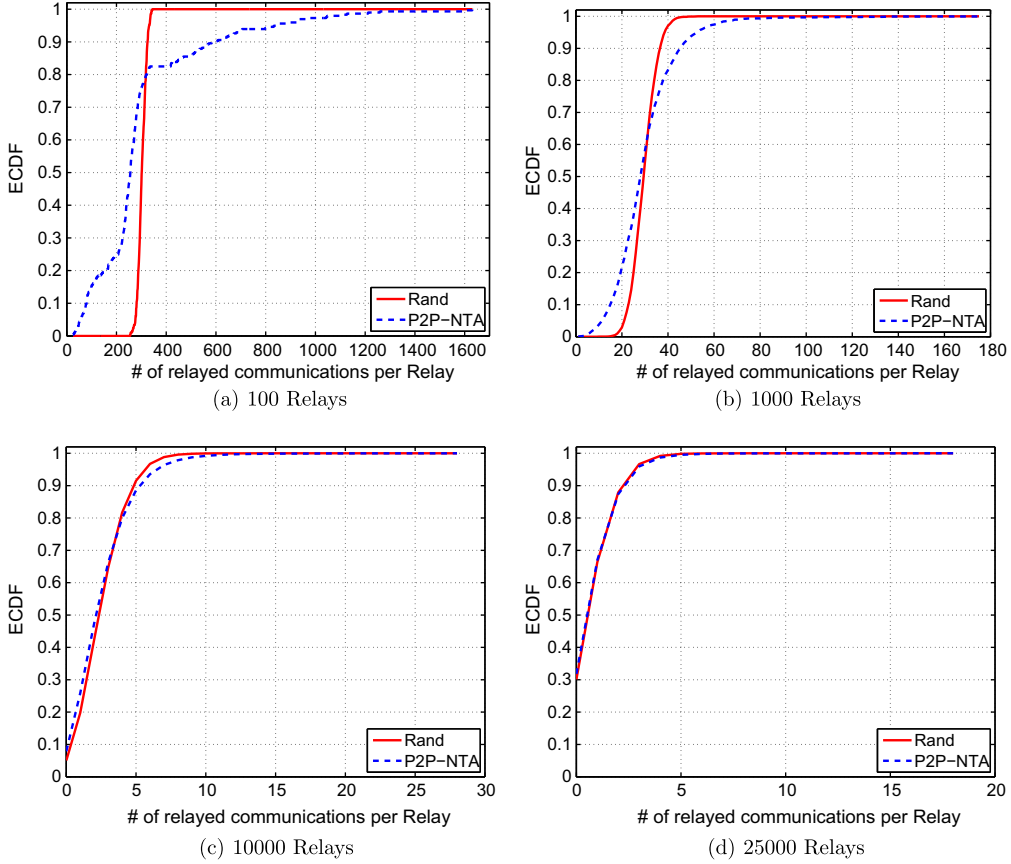


Fig. 9. Number of relayed communications per Relay for different deployments.

shows the ECDF of the Relay lookup latency for the different deployments. As we can observe, the lookup latency is in the order of a few seconds and decreases as we augment the number of Relays. Although this could seem a high value, it is not affecting the QoS of the communications. Note that the P2P-NTA launches the lookup procedure when the application (e.g. Skype) starts. Thus, when the user desires to establish communication the Relay has been already selected. Then, this does not add any extra delay to the standard connection establishment protocol (e.g., [7]) in our solution.

- *Number of supported communications per Relay:* We have computed the number of communications supported by each Relay, for both solutions, P2P-NTA and random Relay selection, considering the different deployments. The resultant ECDFs are presented in Fig. 9. First, we observe that the load produced by our system is similar to that generated by a random Relay selection algorithm. Since a random selection is expected to produce a fair distribution, we can claim that the P2P-NTA is a fair solution in terms of Relay usage. Actually, for large deployments (25,000) the curves are completely overlapped. Hence, the higher the deployment, the fairer the solution is, while still being lightweight.

7. Conclusions

Our work starts from the premise that relayed communications cannot be avoided in today's Internet. Based on real measurements we have shown that using a topologically closer relay has a critical impact on the QoS of the communications and the costs of ISPs. In order to reduce this impact we have introduced: (i) the Gradual Proximity Algorithm (GPA), which finds a topologically close-by available Relay and (ii) the Peer-to-Peer NAT-Traversal Architecture (P2P-NTA), which is a lightweight distributed architecture – based on a DHT – that implements the GPA. We have carried out large-scale simulations using the real Internet topology associated with real delays. The obtained results show that our proposal largely outperforms previous works in the literature. Furthermore, the P2P-NTN exhibits performance levels comparable to direct communication when used by a reasonably large deployment of a p2p application.

Acknowledgement

This work has been partially funded by the Grants MEDIANET (S2009/TIC-1466) from the Regional Government of Madrid and CON-PARTE (TEC2007-67966-C03-03) by the Ministry of Science and Innovation of Spain.

References

- [1] I. Van Beijnum, IPv4 address consumption, Internet Protocol J. 10 (3) (2007) 22–28.
- [2] K. Egevang, P. Francis, RFC 1631 The IP Network Address Translator (NAT), Internet Engineering Task Force, 1994.
- [3] M. Cadaco, M. Freedman, Illuminati – Opportunistic Network and Web Measurement, 2007. <<http://illuminati.coralcdn.org/stats/>>.
- [4] The Impact of P2P File Sharing, Voice over IP, Skype, Joost, Instant Messaging, One-Click Hosting and Media Streaming such as YouTube on the Internet, 2007. <<http://www.ipoque.com/resources/internet-studies/internet-study-2007>>.
- [5] J. Rosenberg, J. Weinberger, C. Huitema, R. Mahy, STUN – Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs), Internet Engineering Task Force, 2003.
- [6] J. Rosenberg, R. Mahy, P. Matthews, Traversal Using Relays Around NAT (TURN): Relay Extensions to Session Traversal Utilities for NAT (STUN), Internet Engineering Task Force, 2009.
- [7] J. Rosenberg, Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols, Internet Engineering Task Force, 2007.
- [8] B. Ford, P. Srisuresh, D. Kegel, Peer-to-peer communication across network address translators, in: ATEC '05: Proceedings of the Annual Conference on USENIX Annual Technical Conference, USENIX Association, 2005, pp. 13–23.
- [9] S. Guha, P. Francis, Characterization and measurement of TCP traversal through NATs and firewalls, in: IMC '05: Proceedings of the 5th ACM SIGCOMM Conference on Internet Measurement, USENIX Association, 2005, p. 18.
- [10] Y. Huang, T.Z. Fu, D.-M. Chiu, J.C. Lui, C. Huang, Challenges, design and analysis of a large-scale p2p-vod system, in: ACM SIGCOMM 2008 Conference on Data Communication, ACM, 2008, pp. 375–388.
- [11] Y. Huang, T.Z. Fu, D.-M. Chiu, J.C. Lui, C. Huang, Challenges, design and analysis of a large-scale p2p-vod system, in: ACM SIGCOMM 2008 Conference on Data Communication, ACM, 2008, pp. 375–388.
- [12] H.V. Madhyastha, T. Isdal, M. Piatek, C. Dixon, T. Anderson, A. Krishnamurthy, A. Venkataramani, iPlane: an information plane for distributed services, in: Proceedings of the 7th Symposium on Operating Systems Design and Implementation, USENIX Association, 2006, pp. 367–380.
- [13] iPlane, iPlane: An Information Plane for Distributed Services. <<http://iplane.cs.washington.edu/>>.
- [14] P. Francis, S. Jamin, C. Jin, Y. Jin, D. Raz, Y. Shavitt, L. Zhang, Idmaps: a global internet host distance estimation service, IEEE/ACM Trans. Networking (2001).
- [15] F. Dabek, R. Cox, F. Kaashoek, R. Morris, Vivaldi: a decentralized network coordinate system, in: ACM SIGCOMM '04, 2004.
- [16] J. Ledlie, P. Gardner, M.I. Seltzer, Network coordinates in the wild, in: USENIX NSDI, 2007.
- [17] S. Agarwal, J.R. Lorch, Matchmaking for online games and other latency-sensitive p2p systems, in: ACM SIGCOMM'09, 2009.
- [18] Y. Rekhter, T. Li, S. Hares, A Border Gateway Protocol 4 (BGP-4), Internet Engineering Task Force, 2006.
- [19] D. Andersen, H. Balakrishnan, F. Kaashoek, R. Morris, Resilient overlay networks, SIGOPS Oper. Syst. Rev. 35 (5) (2001) 131–145.
- [20] K.P. Gummadi, H.V. Madhyastha, S.D. Gribble, H.M. Levy, D. Wetherall, Improving the reliability of internet paths with one-hop source routing, in: OSDI'04: Proceedings of the 6th Conference on Symposium on Operating Systems Design and Implementation, USENIX Association, 2004.
- [21] S. Ren, L. Guo, X. Zhang, Asap: an as-aware peer-relay protocol for high quality voip, in: ICDCS '06: Proceedings of the 26th IEEE International Conference on Distributed Computing Systems, IEEE Computer Society, 2006.
- [22] J.D. Guyton, M.F. Schwartz, in: Locating nearby copies of replicated internet servers, in: ACM SIGCOMM'95, 1995.
- [23] R.L. Carter, M.E. Crovella, Server selection using dynamic path characterization in wide-area networks, in: IEEE INFOCOM '97, 1997.
- [24] D. Katabi, J. Wroclawski, A framework for scalable global ip-anycast (gia), in: ACM SIGCOMM'00, 2000.
- [25] B. Wong, A. Slivkins, E.G. Sirer, Meridian: a lightweight network location service without virtual coordinates, in: ACM SIGCOMM'05, 2005.
- [26] M.J. Freedman, K. Lakshminarayanan, D. Mazières, Oasis: anycast for any service, in: USENIX NSDI'06, 2006.
- [27] P2PSIP Working Group. <<http://www.ietf.org/html.charters/p2psip-charter.html>>.
- [28] J. XingFeng, A Mechanism to Discover STUN/TURN Nodes in P2PSIP, draft-jiang-p2psip-stun-turn-discovery-00, Internet Engineering Task Force, 2007.
- [29] X. Wang, Q. Deng, VIP: a P2P communication platform for NAT traversal, in parallel and distributed processing and applications, in: Third International Symposium, ISPA 2005, 2005, pp. 1001–1011.
- [30] S.A. Baset, H.G. Schulzrinne, An analysis of the skype peer-to-peer internet telephony protocol, in: Proceedings INFOCOM 2006, 25th IEEE International Conference on Computer Communications, IEEE, 2006, pp. 1–11.
- [31] K. Wookyun, S. Baset, H. Schulzrinne, Skype relay calls: measurements and experiments, in: Proceedings INFOCOM Workshops 2008, 27th IEEE International Conference on Computer Communications, IEEE, 2008, pp. 1–6.
- [32] M. Dischinger, A. Mislove, A. Haeberlen, K.P. Gummadi, Detecting bittorrent blocking, in: IMC '08: Proceedings of the 8th ACM SIGCOMM Conference on Internet Measurement, ACM, 2008, pp. 3–8.
- [33] Vuze Network Status Monitor, Azureus. <http://azureus.sourceforge.net/plugin_details.php?plugin=aznetmon>.
- [34] H. Xie, R.Y. Yang, A. Krishnamurthy, Y.G. Liu, A. Silberschatz, P4p: provider portal for applications, SIGCOMM Comput. Commun. Rev. 38 (4) (2008) 351–362.
- [35] D.R. Choffnes, F.E. Bustamante, Taming the torrent: a practical approach to reducing cross-isp traffic in peer-to-peer systems, SIGCOMM Comput. Commun. Rev. 38 (4) (2008) 363–374.
- [36] One-Way Transmission Time, 1996, ITU Rec. G.114, International Telecommunications Union.
- [37] Enterprise QoS Solution Reference Network Design Guide, CISCO, 2005.
- [38] W.B. Norton, The Evolution of the US Internet Peering Ecosystem, Equinix White Papers, 2004.
- [39] S. Kaune, R.C. Rumin, G. Tyson, A. Mauthe, C. Guerrero, R. Steinmetz, Unraveling Bittorrent's File Unavailability: Measurements, Analysis and Solution Exploration, 2009. <<http://arxiv.org/abs/0912.0625>>.
- [40] WHOIS Service. <<http://www.netdemon.net/tutorials/whois.txt>>.
- [41] MaxMind. <<http://www.maxmind.com/>>.
- [42] G. Sigamos, J.M. Pujol, P. Rodriguez, Monitoring the bittorrent monitors: a bird's eye view, in: PAM'09: Proceedings of the Passive Active Measurement Conference, Ser., Lecture Notes in Computer Science, vol. 5448, Springer, 2009, pp. 175–184.
- [43] P.B. Godfrey, I. Stoica, Heterogeneity and load balance in distributed hash tables, in: Proc. of IEEE INFOCOM, 2005.
- [44] S. Surana, B. Godfrey, K. Lakshminarayanan, R. Karp, I. Stoica, Load balancing in dynamic structured peer-to-peer systems, Perform. Eval. 63 (3) (2006) 217–240.
- [45] D.R. Karger, M. Ruhl, Simple efficient load balancing algorithms for peer-to-peer systems, in: SPAA '04: Proceedings of the Sixteenth Annual ACM Symposium on Parallelism in Algorithms and Architectures, ACM, 2004, pp. 36–43.
- [46] R. Cuevas, M. Uruñia, A. Banchs, Routing fairness in chord: analysis and enhancement, in: Proceedings INFOCOM 2009, 28th IEEE International Conference on Computer Communications, IEEE, 2009.
- [47] I. Stoica, R. Morris, D. Liben-Nowell, D.R. Karger, M.F. Kaashoek, F. Dabek, H. Balakrishnan, Chord: a scalable peer-to-peer lookup protocol for internet applications, IEEE/ACM Trans. Network. 11 (1) (2003) 17–32.
- [48] Ookla's Speedtest Throughput Measures. <<http://confluence.slac.stanford.edu/display/IEPM/Ookla%27s+Speedtest+Throughput+Measures>>.
- [49] A. Cabellos-Aparicio, D. Saucz, O. Bonaventure, J. Domingo-Pascual, Validation of a LISP Simulator, UPC, Tech. Rep. UPC-DAC-RR-CBA-2009-8, 2009.



Rubén Cuevas Rumín got his M.Sc. in Telecommunication Engineering from Universidad Carlos III de Madrid in 2005 and his M.Sc. in Network Planning and Management from Aalborg University in 2006. Furthermore, he received an M.Sc. and Ph.D. degree in Telematic Engineering from Universidad Carlos III de Madrid in 2007 and 2010, respectively. Currently, he is Teaching Assistant at the Telematic Engineering Department at Universidad Carlos III de Madrid. From September 2008 until March 2009 he was intern in the Internet Scientific Group at Telefonica Research. His research interests include Overlay and P2P Networks, Content Distribution, Internet Measurements and Online Social Networks.



Ángel Cuevas Rumín got his M.Sc. in Telecommunication Engineering and M.Sc. in Telematic Engineering at Universidad Carlos III de Madrid in 2006 and 2007, respectively. He got an Erasmus Scholarship and completed his Master Thesis at The University of Reading. Currently he is Ph.D. Candidate at the Department of Telematic Engineering at University Carlos III de Madrid. Also, he got a research Internship at SAP Labs France. His research interests include Wireless Sensor Networks, Overlay and P2P Networks and

Optical Networks.



Albert Cabellos-Aparicio received a B.Sc. (2001), M.Sc. (2005) and Ph.D. (2008) degree in Computer Science Engineering from the Technical University of Catalonia (www.upc.edu). In 2002 he joined the Advanced Broadband Communications Center (CCABA, <http://www.ccaba.upc.edu>) where he worked as research assistant. In 2004 he was awarded with a full scholarship to carry out Ph.D. studies at the Department of Computer Architecture, Technical University of Catalonia (UPC), Spain. In September 2005 he

became an Assistant Professor of the Computer Architecture Department and in 2010 he joined the NaNoNetworking Center in Catalunya (<http://www.n3cat.upc.edu>). His main research interests are new architectures for the Internet, nano-networks and wireless measurements. He has participated on several European and national funded research projects.



Loránd Jakab received a degree in Telecommunications Engineering from the Technical University of Cluj-Napoca, Romania in 2004, and is a Ph.D. candidate at the Computer Science Department of the Technical University of Catalonia (UPC) since 2005. His research interests fall into the areas of network mobility and future Internet architectures, with particular focus on the scalability issues of the global inter-domain routing system.



Carmen Guerrero received the Telecommunication Engineering degree in 1994 from the Technical University of Madrid (UPM), Spain, and the Ph.D. in Computer Science in 1998 from the Universidade da Coruña (UDC), Spain. She has been an Assistant Professor (1994–2000) and Associate Professor (2000–2003) at UDC. She is currently Associate Professor since 2003 at Universidad Carlos III de Madrid (UC3M), teaching computer networks courses. She has been involved in several national and international research projects

related with green networking, future Internet, content distribution, overlay networks, information retrieval, broadband access networks, network management and advanced network and multimedia real-time systems. Some of the recent research projects in which she has participated are: CONTENT: Content Home Network and Services for Home Users, MUSE: Multiservice Access Everywhere and E-NEXT: Network of Excellence in Emerging Networking Technologies. More details of her research profile are available at: www.it.uc3m.es/carmen.